

Phylogenetic Tree Computation Tutorial

Frank Olken

Lawrence Berkeley National Lab

Presentation to PGA Course

May 3, 2002

Berkeley, California

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

1

Overview

- Introduction: What? Why? ...
- Multiple Sequence Alignment
- Computing Phylogenetic Trees
- Merging Trees
- Resources
- U.S. Funding Sources

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

2

Introduction

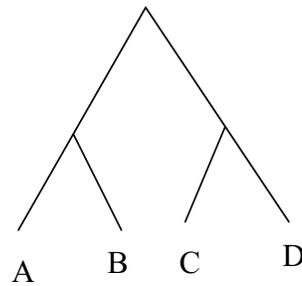
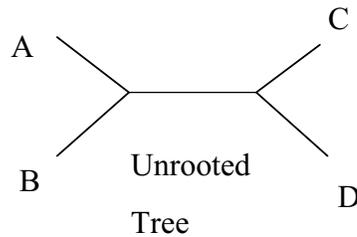
5/6/02

Frank Olken - PGA Phylogeny
Tutorial

3

Example

- Seq. A = A A C C G G T T
- Seq. B = A A C C G G T G
- Seq. C = A C C C G G T C
- Seq. D = A C C C G G T A



5/6/02

Frank Olken - PGA Phylogeny
Tutorial

4

Similarity vs. Homology

- Similar
 - sequences resemble one another
- Homolog
 - sequences derived from common ancestor
- Ortholog
 - homologous sequences within a species
- Paralog
 - homologous sequences between species

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

5

Ortholog vs. Paralog

- Ortholog – genomic variation occurs after speciation – hence can be used for phylogeny of organism
- Paralog – genetic duplication occurs before speciation – hence not suitable for phylogeny of organism

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

6

Homoplasy

- Sequence similarity NOT due to common ancestry
- May arise due to parallelism or convergent evolution
- Parallelism
- Convergent evolution

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

7

Phylogenetic Tree

- Binary tree (i.e, fan-out from nodes = 2)
- Variously rooted or unrooted
- Tree topology
- Branch lengths (evolutionary time)

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

8

Phylogeny of what?

- **Organisms**
 - Whole genome phylogeny
 - Ribosomal RNA (surrogate for whole genome)
- **Strains (closely related microbes)**
- **Individual genes (or gene families)**
- **Repetitive DNA sequences**
- **Metabolic pathways**
- **Secondary Structures**
- **Any discrete character(s)**
- **Human languages**
- **Microbial communities**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

9

Why compute phylogenetic trees?

- **Understand evolutionary history**
- **Map pathogen strain diversity for vaccines**
- **Assist in epidemiology**
 - Of infectious diseases
 - Of genetic defects
- **Aid in prediction of function of novel genes**
- **Biodiversity studies**
- **Understanding microbial ecologies**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

10

Rooted vs. Unrooted Tree

- Root - ancestor of all taxa considered
- Unrooted tree - typical result, unknown common ancestor
- Rooted tree - known common ancestor
- Specify root by means of outgroup
- Outgroup is distant from all other taxa
 - example: mammals and a salamander
 - ancestor of outgroup is presumed root

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

11

Which Sequences ?

- DNA
 - Very sensitive, non-uniform mutation rates
- cDNA/RNA
 - Useful for more remote homologies
- Protein Sequences
 - Useful for most remote homologies, deep phylogenies, more uniform mutation rates, more character states

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

12

Ribosomal RNA 16S Sequences

- These sequences exist in all organisms
- They are highly conserved
- Hence suitably for broad, very deep phylogeny studies
- Compiled for tens of thousands of organisms, mostly microbial
- Unsuitable to fine grained phylogeny

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

13

What is to be computed?

- Tree topology
 - branching order
 - root
- Branch lengths (evolutionary time)
- Ancestral sequences
- Tree figure of merit (e.g., likelihood)
- Tree reliability

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

14

Computational Process

- Get DNA/RNA/Protein Sequences
- Construct multiple sequence alignment
- Compute pairwise distances
 - (for distance methods)
- Build tree: topology + branch lengths
- Estimate reliability
- Visualize

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

15

Multiple Sequence Alignment

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

16

Multiple Sequence Alignment

- Align sequences so that corresponding positions are in same columns
- Pad missing nucleotides as nulls where needed
- Each “column” then becomes a single character in the phylogeny computation

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

17

Multiple Sequence Alignment

- NP-hard
 - Worst case computational complexity is exponential in number of problem size
- In practice, often use greedy algorithms
 - choose best incremental change in solution
 - no backtracking
- Figure of merit:
 - Sum of pair-wise edit distances

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

18

Multiple Sequence Alignment Algorithms

- Dynamic Programming
- Hidden Markov Models
- Stochastic Context Free Grammars

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

19

Dynamic Programming MSA - Sequential Addition

- Compute pairwise edit distances between sequences via dynamic programming
- Merge closest pair
- Generate consensus sequence
- Merge sequence closest to consensus sequence alignment

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

20

Dynamic Programming MSA - Bottom Up Merging

- Compute pairwise edit distances between sequences via dynamic programming
- Merge closest pair
- Replace merged pair with consensus sequence
- Recompute pairwise edit distances, and merge closest pair

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

21

Hidden Markov Model MSA

- Estimate HMM for set of sequences
- Compute most probable alignment of each sequence to HMM
- Use this as basis for MSA
- HMM allows probabilistic representation of “consensus sequence”

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

22

SCFG MSA Algorithm

- Estimate Stochastic Context Free Grammar for set of ribosomal RNA sequences (e.g., 16S sequences)
- Compute most probable sequence alignment to SCFG for each sequence
- Use this as basis of multiple sequence alignment
- Preserves RNA secondary structure in MSA

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

23

Algorithms for computing phylogenetic trees

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

24

Evaluation Criteria for Tree Computational Methods

- Accuracy
- Explicit statistical model of evolution ?
- Efficient use of data
- Computation Time
- Branch lengths ?
- Quality measure ?
- Reliability measure ?

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

25

Computational Approaches to Phylogenetic Tree Computation

- Distance-based methods
 - UPGMA, Neighbor joining
- Maximum Parsimony Method
- Maximum Likelihood Methods
- Tree merging
 - Consensus trees, supertrees

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

26

Distance vs. Character State Methods

- Distance Methods
 - UPGMA, Neighbor Joining, Min. Evol.,
 - Requires distance measures between sequences
 - Suitable for continuous characters
- Character State Methods
 - Max. parsimony, Max. Likelihood, ...
 - Requires discrete characters

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

27

Similarity Measures vs. Distances

- Similarity Measure
 - bigger value = more similar
- Distances
 - bigger value = less similar
 - triangle inequality
 - $|x,y| + |y,z| < \text{or} = |x,z|$
 - often assumed additive for distance-based phylogenetic tree construction

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

28

Can distance-based methods use similarity measures?

- Maybe ...
- Depends on whether distance methods uses:
 - triangle inequality
 - additive distance measure

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

29

Simple distances between sequences

- Number of differing positions
- Weighted differences
- Edit distances (weighted sum of inserts, deletes, substitutions)
- Weighted Substitution Cost Matrices
 - PAM, BLOSUM
- Poisson Corrections (next slide)

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

30

Distance Metric Between Sequences

- $p = n_d / n$
- = number of characters which differ / total number of characters
- p is not proportional to evolutionary time
- Reason: sites can mutate more than once
- Poisson correction:
- $d = -\ln(1-p)$

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

31

Similarity Measure for Protein Structures

- **Construct contact map (graph) for each protein structure**
 - vertex = residue (amino acid)
 - edge = distance between AA's is less than 5 Angstroms
- **Compute pairwise alignment between structures**
 - nonoverlapping matching of residues
- **Similarity Measure**
 - = number of shared edges from contact graphs

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

32

UPGMA Distance Method

- Unweighted Pair Group Method Using Arithmetic Mean
- by Sokal and Michener, 1958
- Merge closest pair of taxa (by distance)
- Recompute distances to merged node via mean of pairwise distances to leaves
- Repeat

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

33

UPGMA Method

- Fast to compute
- Implicitly assume molecular clock
 - i.e., uniform mutation rates across sites and branches

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

34

Neighbor Joining Distance Method

- Compute pairwise distances, $d(i,j)$, set $L =$ all leaves T
- Compute $D(i,j) = d(i,j) - (r(i)+r(j))$
- $r(i) =$ average distance to other leaves
- Merge closest pair of sequences i and j
 - for new k , set $d(k,m) = 1/2 (d(i,m)+d(j,m)-d(i,j))$ for m in L
 - Add k to T with
 - set $d(i,k) = 1/2 (d(i,j)+r(i)-r(j))$
 - set $d(j,k) = d(i,j)-d(i,k)$
 - replace i and j with k in L
- Repeat

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

35

Neighbor Joining

- Generates unrooted trees
- Assumes additive distances in tree

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

36

Merits of Distance Methods

- Fastest method
- Not very accurate or efficient use of data
- Can use continuous data (not just sequences)
- No statistical model of evolution
- No figure of merit
- No branch lengths

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

37

Maximum Parsimony

- Minimize total number of state changes along all paths in tree
- Intermediate computational cost
- Figure of merit = number of state changes
- No statistical model
- No branch lengths

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

38

Maximum Likelihood

- $L = p(\text{data} \mid \text{tree, branch lengths, model})$
- Search trees, branch lengths to find max L
- This is MLE tree
- Algorithms:
 - generate tree topology
 - optimize branch lengths
 - retain if best seen
 - loop

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

39

Maximum Likelihood Methods

- Explicit statistical model
- Figure of merit = log likelihood
- Computes branch lengths
- Very expensive to compute, use heuristics
- Reliability estimation by resampling
- Efficient use of sequence data
- Examples: Phylip/dnaML, fastdnaML

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

40

Maximum Likelihood Estimation - Assumptions

- Characters (nucleotide positions) evolve independently
- Mutation Rate variation:
 - Molecular clock \implies uniform rates across positions and branches
 - We can allow rate to vary by position (usually assume Gamma distribution)
 - Requires that estimate more parameters

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

41

MLE Phylogenetic Tree Construction Algorithms

- Must compute both tree topology, branch lengths, rates of evolution (for each position)
- Too many tree topologies to search exhaustively
- Hence heuristic tree search
 - Order taxa (randomly?)
 - Build tree incrementally
 - For each taxa, consider all possible locations in existing tree to insert
 - Compute likelihood for each possible insertion point.
 - Choose best insertion point
 - Get next taxa

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

42

MLE Phylogenetic Tree Algorithms

- Branch swapping
 - Given tree
 - Consider all possible pairwise branch swaps, where branches within distance “k” in tree
 - For each possible swap
 - reoptimize branch lengths
 - compute likelihood of resulting tree

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

43

MLE Phylogenetic Tree Algorithms

- Overview of algorithm
 - Pick next taxa
 - Generate possible topologies
 - Perform branch optimization for each topology
 - Compute likelihood of resulting tree
 - Retain best “k” trees
 - Get next taxa

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

44

Popular MLE Codes

- dnaML - Joe Felsenstein (U. Washington)
- fastdnaML - Gary Olsen (UIUC)
- PAUP - Dave Swofford (Florida State U.)
- PAML

Merging Phylogenetic Trees

Merging Phylogenetic Trees

- Consensus Trees
 - All trees computed over same taxa
- Supertrees
 - Trees computed over overlapping sets of taxa

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

47

Consensus Trees

- All trees computed over same set of taxa
- Used for cases when there exist many tree with similar figures of merit (parsimony scores, likelihood estimates)
- Goal: identify robust result

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

48

Consensus Trees

- **Strict Consensus**
 - Only report branch orders which occur in all trees
 - Use multi-way branching when branch order is not consistent across all trees
 - Very strong consensus criteria

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

49

Supertrees

- Supertrees are constructed from candidate trees computed over overlapping sets of taxa
- Supertree construction is a way to compute very large trees.
- Usually need significant pairwise overlaps among at least some pairs of trees taxa

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

50

Majority Consensus Trees

- Use voting algorithms on branch orders
- Report branch orders which supported by a majority (or supermajority, e.g., 2/3) of candidate trees
- More relaxed than strict consensus trees
- Widely used

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

51

Supertrees

- Merge overlapping phylogenetic trees
- DCM Algorithm
 - Tandy Warnow at UT Austin

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

52

Tree Reliability

- Bootstrapping (Felsenstein, et al.)
 - Resample characters (nucleotide positions)
 - Recompute phylogenetic tree
 - Repeat
 - Compare trees, compute consensus tree
 - Expensive
 - Sample size = 100 (?)

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

53

Whole Genome Phylogeny

- Done at gene level (not sequence)
- Look at larger genome rearrangements
- Good for deep phylogeny
- Consider inversions, translocations
- Currently primarily suitable for microbes
 - most whole genomes sequences have been for microbes
- See papers by Sankhoff, et al.

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

54

Choosing an algorithm

- **Continuous characters, lots of data, few computing resources ==> Neighbor joining**
- **Discrete characters, few mutations / homoplasy ==> Maximum Parsimony**
- **Discrete characters, limited sequence lengths, some homoplasy ==> Maximum Likelihood Estimation**
- **Discrete characters, many taxa ==> Supertree**
- **Complete genomes ==> Whole genome phylogeny**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

55

Choosing an algorithm

- **For low numbers of mutations - MP and MLE generate similar results**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

56

Practical problems

- Often difficult to resolve deep branch orders
 - (ancient branching of evolutionary tree)
- Tradeoff between:
 - Sensitivity (rapid genome evolution)
 - Detection of ancient branching (slow genome evolution)

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

57

Data Management for Phylogeny

- Sequence Data Management
- Multiple Sequence Alignments
- Tree Data Management
- Tree Alignments
 - consensus subtrees

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

58

Phylogeny for Vaccine Studies

- Phylogeny of disease strains
- Need to assure that vaccine confers resistance to all strains
- Phylogeny is used to select diverse strains for constructing vaccine
- Example: influenza

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

59

Phylogeny and Epidemiology

- Pathogen phylogeny used to assist epidemiological studies
- Example: HIV
 - rapid evolution of virus
 - use phylogeny to verify source of infection of particular individual
- Co-evolution of pathogens and hosts
- See Crandall, *Evolution of HIV*

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

60

Resources on Phylogeny

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

61

Web Resources

- **Felsenstein's Phylogenetic Program Directory**
 - <http://evolution.genetics.washington.edu/phylip.html>
- **UT Austin Phylogenetics Lab**
 - <http://kristin.csres.utexas.edu/>
- **Woese Lab**
 - <http://www.life.uiuc.edu/micro/woese.html>
- **Tree-of-life web site**
 - <http://tolweb.org/tree/phylogeny.html>

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

62

Bibliographies and Search Engines

- **Reading lists to IB 200A + IB 200B**
 - by Mishler at UC Berkeley
 - <http://ib.berkeley.edu/courses/ib200a/>
- **Medline (a.k.a. PubMed) at NLM**
- **BIOSIS (Biological Abstracts)**
- **INSPEC (by IEE)**
- **NEC Citeseer (CS papers)**
- **MathSciNet (Mathematical Reviews)**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

63

Books on Phylogeny

- Graur, Li. *Fundamentals of Molecular Evolution*, Sinauer
- Hall, *Phylogenetics Made Easy*, Sinauer
- Hillis, Moritz, Mable. *Molecular Systematics*, 2nd edition, Sinauer, 1996
- Kitching, Forey, Humphries. *Cladistics: The Theory and Practice of Parsimony Analysis*, 1998
- Kimura, M. *The Neutral Theory of Evolution*, Cambridge, 1983
- Li. *Molecular Evolution*, Sinauer
- Nei, M. & S. Kumar. *Molecular Evolution and Phylogenetics*, Oxford, 2000
- Page & Holmes. *Molecular Evolution: A Phylogenetic Approach*, 1998
- Smith, J.M., *Evolutionary Genetics*, 1998
- Wheeler & Meier. *Species Concepts and Phylogenetic Theory*, 2000
- Wilkins. *Evolution of Developmental Pathways*, Sinauer, 2001

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

64

More Books on Phylogeny

- Harvey, Leigh Brown, Smith, Nee. *New Uses for New Phylogenies*, Oxford, 1966
- Crandall, K. (editor) *The Evolution of HIV*, Johns Hopkins Univ. Press, 1999
- Mount, D.W. *Bioinformatics: Sequence and Genome Analysis*, Cold Spring Harbor Laboratory Press, 2000, Chapter 6 - Phylogenetic Prediction
- Doolittle, R.F. *Computer Methods for Macromolecular Sequence Analysis*, Methods in Enzymology, vol. 266, 1996, Academic Press

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

65

Journals on Phylogeny

- **Cladistics**
- **Molecular Biology and Evolution**
- **Molecular Phylogenetics and Evolution**
- **Systematic Biology**
- **Systematic Zoology**
- **Evolutionary Biology**
- **Taxon**
- **Bioinformatics**
- **J. of Computational Biology**
- **J. of Theoretical Biology**
- **Ecology and Evolutionary Biology**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

66

Conferences on Phylogeny

- **RECOMB**
- **ISMB (Intelligent Systems for Molecular Biology)**
- **Evolution 2002**
- **Classification Society of N. America Annual Mtg.**
- **Conf. Of the Int'l. Federation of Classification Societies (IFCS)**
- **ICSEB (Intl. Conf. On Systematics & Evolutionary Biology)**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

67

Phylogenetics Organizations

- **<http://www.ucmp.berkeley.edu/subway/phylo/phyloorg.html>**
- **Classification Society of North America**
- **Society for Molecular Biology and Evolution**
- **Will Hennig Society**
- **Classification Society**
- **Green Plant Phylogeny Group**
- **European Society for Evolutionary Biology**
- **Society for Systematic Biology**
- **Systematics Society (Britain)**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

68

Funding for Phylogeny Research

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

69

U.S. Funding Sources for Computational Phylogeny

- NSF
- NIH
- CDC
- DOE
- DOD

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

70

Funding - NSF

- **NSF Biology Directorate**
 - Primary funder for phylogeny in U.S.
 - Both specific programs in phylogeny
 - and **ITR (Info. Tech. Research) Program** (large interdisciplinary projects)
 - Contact Sylvia Spengler
 - **Tree of Life Program** - Diana Lipscomb

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

71

Funding - NIH

- **National Institutes of Health**
- **Primarily interested in phylogeny and human health**
- **Phylogeny of pathogens, disease genes, ...**
- **NIAID - National Institute of Allergies and Infectious Diseases**
 - **phylogeny of pathogens of infectious diseases**
- **NIGMS - National Institute of General Medical Sciences**
 - **phylogeny of gene families**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

72

Funding - CDC

- Center for Disease Control (Atlanta)
- Concerned with infectious diseases and food toxins
- Phylogeny of pathogens
- Phylogeny and epidemiology

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

73

Funding - DOE

- U.S. Department of Energy
- Microbial Research Program
 - Prokaryotes primarily
 - Bioremediation
 - Carbon sequestration
 - Biofuels
- Microbial Phylogeny & Ecology
- Bioterrorism / Biowarfare Program (NN20)

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

74

Funding - DOD

- Biowarfare / Bioterrorism concerns
- Phylogeny of pathogens
- Understanding Pathogenicity
- DARPA Biowarfare Defense
 - **DARPA DSO Pathogen Genomic Sequencing (Dr. Eric Eisenstadt)**

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

75

Acknowledgements

- This research supported by U.S. Department of Energy, Office of Biological and Environmental Research

5/6/02

Frank Olken - PGA Phylogeny
Tutorial

76

Contact Information

- **Frank Olken**
 - **Lawrence Berkeley National Laboratory**
 - **1 Cyclotron Road, Mailstop 50B-3238**
 - **Berkeley, CA 94720**
 - **Tel: 510-486-5891**
 - **Fax: 510-586-4004**
 - **Email: olken@lbl.gov**
 - **URL: <http://www.lbl.gov/~olken>**